

# The Benefits of Network Coding for Peer-to-Peer Storage Systems

Alexandros G. Dimakis, P. Brighten Godfrey, Martin J. Wainwright and Kannan Ramchandran

Department of Electrical Engineering and Computer Science,  
University of California, Berkeley, CA 94704.

Email: {adim, pbg, wainwrig, kannanr}@eecs.berkeley.edu

**Abstract**—Peer-to-peer distributed storage systems provide reliable access to data through redundancy spread over nodes across the Internet. A key goal is to minimize the amount of bandwidth used to maintain that redundancy. Storing a file using an erasure code, in fragments spread across nodes, promises to require less redundancy and hence less maintenance bandwidth than simple replication to provide the same level of reliability. However, since fragments must be periodically replaced as nodes fail, a key question is how to generate a new fragment in a distributed way while transferring as little data as possible across the network.

In this paper, we introduce a general technique to analyze storage architectures that combine any form of coding and replication, as well as presenting two new schemes for maintaining redundancy using network coding. First, we show how to optimally generate MDS fragments directly from existing fragments in the system. Second, we introduce a new scheme called Regenerating Codes which use slightly larger fragments than MDS but have lower overall bandwidth use. We also show through simulation that in realistic environments, Regenerating Codes can reduce maintenance bandwidth use by 25% or more compared with the best previous design—a hybrid of replication and erasure codes—while simplifying system architecture.

## I. INTRODUCTION

The purpose of distributed file storage systems such as OceanStore [17], Total Recall [3], and DHash++ [6] is to store data reliably over long periods of time using a distributed collection of disks (say, at various nodes across the Internet). Ensuring reliability requires the introduction of redundancy, the simplest form of which is straightforward replication.

Several designs [16], [3], [6] use erasure codes instead of replication. A *Maximum-Distance Separable* (MDS) erasure code stores a file of size  $M$  bytes in the form of  $n$  fragments each of size  $M/k$  bytes, any  $k$  of which can be used to reconstruct the original file.

However, a complication arises: in distributed storage systems, redundancy must be continually refreshed as nodes choose to leave the system and disks fail, which involves large data transfers across the network. How do we efficiently create new encoded fragments in response to failures? A new replica may simply be copied from any other node storing one, but traditional erasure codes require access to the original data to produce a new encoded fragment. How do we generate an erasure encoded fragment when we only have access to erasure encoded fragments?

In the *naive strategy*, the node which will store the new fragment—which we will call the *newcomer*—downloads  $k$  fragments and reconstructs the file, from which a new fragment is produced. Thus,  $M$  bytes are transferred to generate a fragment of size only  $M/k$ .

To reduce bandwidth use, one can adopt what we call the *Hybrid strategy* [18]: one full replica is maintained in addition to multiple erasure-coded fragments. The node storing the replica can produce new fragments and send them to newcomers, thus transferring just  $M/k$  bytes for a new fragment. However, maintaining an extra replica on one node dilutes the bandwidth-efficiency of erasure codes and complicates system design. For example, if the replica is lost, new fragments cannot be created until it is restored. In fact, one study comparing the Hybrid strategy with replication in distributed storage systems [18] argued that in practical environments, Hybrid’s reduced bandwidth is limited, and may be outweighed by its drawbacks, in part due to the added complication of maintaining two types of redundancy.

It is thus natural to pose the following question: is it possible to maintain an erasure code using less bandwidth than the naive strategy, without resorting to an asymmetric strategy like Hybrid? More deeply, what is the minimal amount of data that must be downloaded in order to maintain an erasure code?

In this paper we show how network coding can help for such distributed storage scenarios. We introduce a general graph-theoretic framework through which we obtain lower bounds on the bandwidth required to maintain any distributed storage architecture and show how random linear network coding can achieve these lower bounds.

More specifically, we determine the minimum amount of data that a newcomer has to download to generate an MDS or nearly-MDS fragment, a scheme which we call *Optimally Maintained MDS* (OMMDS). In particular, we prove that if the newcomer can only connect to  $k$  nodes to download data for its new fragment, then the  $M$ -byte download of the naive strategy is the information-theoretic minimum. Surprisingly, if the newcomer is allowed to connect to more than  $k$  nodes, then the total download requirement can be reduced significantly. For example, if  $k = 7$  (the value used in DHash++ [6]),  $n = 14$ , and a newcomer connects to  $n - 1$  nodes, a new fragment can be generated by transferring  $0.27M$  bytes, or 73% less than the naive strategy. However, the associated overhead is still substantial, and it turns out that Hybrid offers a better reliability-bandwidth tradeoff than OMMDS. To improve on Hybrid, we must therefore look beyond MDS codes.

With this perspective in mind, we introduce our second scheme, *Regenerating Codes* (RC), which minimize amount of data that a newcomer must download subject to the restriction that we preserve the “symmetry” of MDS codes. At a high level, the RC scheme improves on OMMDS by having a newcomer store all the data that it downloads, rather than throwing some away. As a consequence, RC has slightly *larger* fragments than MDS, but very low maintenance bandwidth

overhead, even when newcomers connect to just  $k$  nodes. For example, if  $k = 7$ , a newcomer needs to download only 0.16M bytes—39% less than OMDS and 84% less than the naive strategy. Moreover, our simulation results based on measurements of node availability in real distributed systems show that RC can reduce bandwidth use by up to 25% compared with Hybrid when  $k = 7$ . RC improves even further as  $k$  grows.

We emphasize that there are still tradeoffs between RC and other strategies. For example, users wishing to reconstruct the file pay a small overhead due to RC’s larger fragments. Nevertheless, RC offers a promising alternative due to its simplicity and low maintenance bandwidth.

## II. RELATED WORK

Network coding for distributed storage was introduced in [7], [8] in a sensor network scenario where a bandwidth was minimized for a static setup. Related work for storage in wireless networks includes [10], [14], [24], [21], [1].

Network coding was proposed for peer-to-peer content distribution systems [11] where random linear operations over packets are performed to improve downloading. Random network coding was also recently proposed for P2P network diagnosis [23]. Our paper is based on similar ideas but the storage systems have different performance metrics that need to be analyzed.

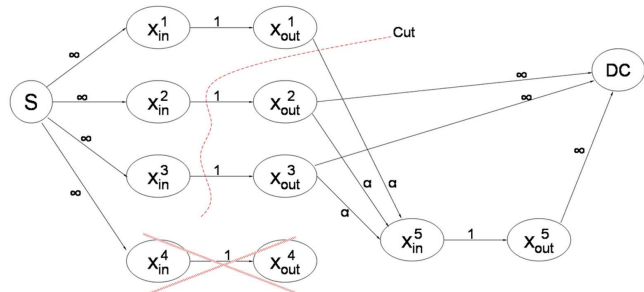
A number of recent studies [22], [3], [18] compared replication with erasure coding for large-scale, peer-to-peer distributed storage systems. The analysis of Weatherspoon and Kubiatowicz [22] showed that erasure codes reduced storage and bandwidth use by an order of magnitude compared with replication. Bhagwan et al [3] came to a similar conclusion in a simulation of the Total Recall storage system. However, Rodrigues and Liskov [18] show that in high-churn (i.e., high rate of node turnover) environments, erasure coding provides a large benefit but the maintenance bandwidth cost is too high to be practical for a P2P distributed storage system. In low-churn environments, the reduction in bandwidth is negligible. In moderate-churn environments, there is some benefit, but this may be outweighed by the added architectural complexity that erasure codes introduce. These results [18] apply to the Hybrid strategy. In Section V, we repeat the evaluation of [18] to measure the performance of the two redundancy maintenance schemes that we introduce.

## III. FUNDAMENTAL LIMITS ON BANDWIDTH

### A. Information flow graph

Our analysis is based on a particular graphical representation of a distributed storage system, which we refer to as an *information flow graph*  $\mathcal{G}$ . This graph describes how the information of the data object travels through time and storage nodes and reaches reconstruction points at the data collectors. More precisely, it is a directed acyclic graph consisting of three kinds of nodes: a single data source  $S$ , storage nodes  $x_{in}^i, x_{out}^i$  and data collectors  $DC_i$ . The single node  $S$  corresponds to the source of the original data. Storage node  $i$  in the system is represented by a storage input node  $x_{in}^i$ , and a storage output node  $x_{out}^i$ ; these two nodes are connected by a directed edge  $x_{in}^i \rightarrow x_{out}^i$  with capacity equal to the amount of data stored at node  $i$ . See Figure III-A for an illustration.

Given the dynamic nature of the storage systems that we consider, the information flow graph also evolves in time. At any given time, each vertex in the graph is either *active* or *inactive*, depending on whether it is available in the network. At the initial time, only the source node  $S$  is active; it then contacts an initial set of storage nodes, and connects to their inputs ( $x_{in}$ ) with directed edges of infinite capacity. From this point onwards, the original source node  $S$  becomes and remains inactive. At the next time step, the initially chosen storage nodes become now active; they represent a distributed erasure code, corresponding to the desired steady state of the system. If a new node  $j$  joins the system, it can only be connected with active nodes. If the newcomer  $j$  chooses to connect with active storage node  $i$ , then we add a directed edge from  $x_{out}^i$  to  $x_{in}^j$ , with capacity equal to the amount of data that the newcomer downloads node  $i$ . Note that in general it is possible for nodes to download more data than they store, as in the example of the (14, 7)-erasure code. If a node leaves the system, it becomes inactive. Finally, a data collector  $DC$  is a node that corresponds to a request to reconstruct the data. Data collectors connect to subsets of active nodes through edges with infinite capacity.



**Fig. 1.** Illustration of an information flow graph  $\mathcal{G}$ . Suppose that a particular distributed storage scheme uses an (4, 3) erasure code in which any 3 fragments suffice to recover the original data. If node  $x^4$  becomes unavailable and a new node joins the system, then we need to construct new encoded fragment in  $x^5$ . To do so, node  $x_{in}^5$  is connected to the  $k = 3$  active storage nodes. Assuming that it downloads  $\alpha$  bits from each active storage node, of interest is the minimum  $\alpha$  required. The min-cut separating the source and the data collector must be larger than 3 for reconstruction to be possible. For this graph, the min-cut value is given by  $2 + \alpha$ , implying that  $\alpha \geq 1$ , so that the newcomer has to download the complete data object if he connects to only  $k = 3$  storage nodes.

### B. Bounds

To obtain bounds on the how much each storage node has to download, we use the following lemma. Due to space constraints we will only present sketches or fully omit some proofs.

*Lemma 1:* A data collector  $DC$  can never reconstruct the initial data object if the minimum cut in  $\mathcal{G}$  between  $S$  and  $DC$  is smaller than the initial object size.

The next claim, which builds on known results from network coding, shows that there exist linear network codes which can match this lower bound for all data collectors, and also that simple linear mixing of packets using random independent co-

efficient over a finite field (randomized network coding [13]) will be sufficient with high probability.

*Proposition 1:* Assume that for some distributed storage scheme, we construct the  $\mathcal{G}$  graph and place all the possible  $\binom{n}{k}$  data collectors where  $n$  is the number of active nodes. If the minimum of the min-cuts separating the source with each data collector is at least the data object size  $\mathcal{M}$ , then there exists a linear network code such that all data collectors can recover the data object. Further, randomized network coding guarantees that all collectors can recover the data object with probability that can be driven arbitrarily high by increasing the field size.

*Proof:* (sketch) This proof is based on a reduction of the distributed storage problem into a multicasting problem with a single source sending its data to all  $\binom{n}{k}$  possible data collectors. We can then apply known results for single-source multicast; network coding can achieve the associated min-cut/max-flow bound [2] and from [15] we know that a linear network code will suffice.

Ho et al. [13] show that the use of random linear network codes at all storage nodes suffices to ensure that each data collector can reconstruct with probability that can be pushed arbitrarily high by increasing the field size. (See in particular Theorem 3 in the paper [13], which ensures that the probability is at least  $(1 - \frac{d}{q})^N$ , where  $d$  is the number of data collectors and  $N$  is total number of storage nodes in  $\mathcal{G}$  and  $q$  is the field size.) ■

The above results allow us to provide a complete characterization of the bandwidth cost associated with maintaining an MDS erasure code:

*Proposition 2:* Assume the data object is divided in  $k$  fragments, an  $(n, k)$ -MDS code is generated and one encoded fragment is stored at each node. Suppose that one node leaves the system and that a new joining node wishes to create a new encoded fragment by downloading an  $\alpha$  fraction of a fragment from each of  $n - 1$  active storage nodes. Then we must have  $\alpha \geq \frac{1}{n-k}$  for successful reconstruction.

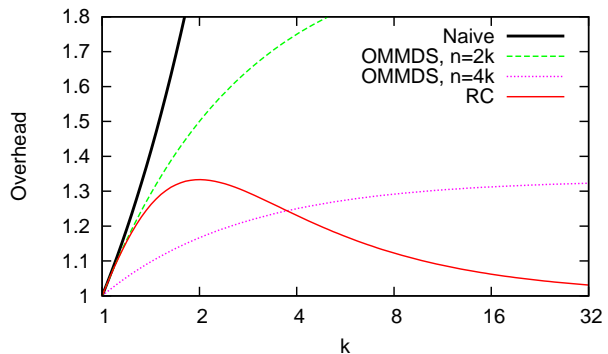
*Proof:* Consider the information flow graph  $\mathcal{G}$  for this storage system. Suppose that any newcomer connects to  $n - 1$  storage nodes and downloads a portion  $\alpha$  of the fragment from each storage node, where  $\alpha$  is to be determined. A data collector is connected to the newcomer and  $k - 1$  other storage nodes. The minimum cut in this newly formed  $\mathcal{G}$  is given by  $k - 1 + (n - 1 - (k - 1))\alpha$ ; for successful reconstruction, it has to be at least  $k$ , so  $\alpha \geq \frac{1}{n-k}$  is the minimum possible bandwidth to maintain an MDS code. ■

In the special case of the  $(n, k) = (14, 7)$  erasure code considered in the introduction, Proposition 2 verifies the earlier claim that the newcomer needs to download only  $\frac{1}{7}$  of a fragment from each of the  $n - 1 = 13$  active storage nodes, for a total of  $\frac{1}{7} \cdot \frac{1}{k} \mathcal{M}(n - 1) \approx 0.27\mathcal{M}$  bytes.

We refer to MDS codes maintained in the procedure specified by Proposition 2 as *Optimally Maintained MDS*, or *OMMDS* for short.

#### IV. REGENERATING CODES

The OMMDS scheme of the previous section is a significant improvement over the naive scheme of downloading the entire file to generate a new fragment. However, the associated overhead is still substantial, and our experimental evaluation in Section V reveals that the Hybrid scheme still offers a better



**Fig. 2.** The *overhead*  $\beta$  is the number of bytes downloaded to produce a fragment, divided by the size of an MDS fragment. For the naive strategy,  $\beta_{\text{naive}} = k$ ; for OMMDS in which newcomers connect to  $n - 1$  nodes,  $\beta_{\text{OMMDS}} = \frac{n-1}{n-k}$ ; for RC in which newcomers connect to just  $k$  nodes,  $\beta_{\text{RC}} = k^2 / (k^2 - k + 1)$ . Moreover, RC fragments are  $\beta_{\text{RC}}$  times larger than MDS fragments, so that the data collector must download  $\beta_{\text{RC}}$  times the size of the original file.

reliability-bandwidth tradeoff than the OMMDS. Moreover, as established in Proposition 2, an MDS code cannot be maintained with less bandwidth than OMMDS. Therefore, we can only hope to use less bandwidth with a coding scheme other than an MDS code.

With this perspective in mind, this section introduces the notion of a *Regenerating Code* (RC). Subject to the restrictions that we preserve the “symmetry” of MDS codes (to be detailed later), we derive matching lower and upper bounds on the minimal amount of data that a newcomer must download. In contrast with OMMDS, the RC approach has very low bandwidth overhead, even when newcomers connect to just  $k$  nodes. At a high level, the RC scheme improves on OMMDS by having a newcomer store *all* the data that it downloads, rather than throwing some away. As a consequence, RC fragments are slightly larger than MDS fragments, by a factor  $\beta_{\text{RC}} = k^2 / (k^2 - k + 1)$  (see Figure 2 for an illustration), and any data collector that reconstructs the file downloads  $\beta_{\text{RC}}$  times the size of the file. However, note that  $\beta_{\text{RC}} \rightarrow 1$  as  $k \rightarrow \infty$ .

Regenerating codes minimize the required bandwidth under a “symmetry” requirement over storage nodes. Specifically, we require that any  $k$  fragments can reconstruct the original file; all fragments have equal size  $\alpha\mathcal{M}$ ; and a newcomer produces a new fragment by connecting to any  $k$  nodes and downloading  $\alpha\mathcal{M}/k$  bits from each. In this paper, to simplify the scheme, we fix the number of nodes to which the newcomer connects to  $k$  (the minimum possible). The free parameter  $\alpha$  will be chosen to minimize bandwidth.

Assume that newcomers arrive sequentially, and that each one connects to an arbitrary  $k$ -subset of previous nodes (including previous newcomers). The following result characterizes the bandwidth requirements of the RC scheme:

*Theorem 1:* Assume all storage nodes store  $\alpha\mathcal{M}$  bits and newcomers connect to  $k$  existing nodes and download  $\frac{1}{k}\alpha\mathcal{M}$  bits from each. Then, define

$$\alpha_c = \frac{1}{k} \times \frac{1}{1 - \frac{1}{k} + \frac{1}{k^2}}. \quad (1)$$

If  $\alpha < \alpha_c$  then reconstruction at some data collector who

connects to  $k$  storage nodes is information theoretically impossible.

If  $\alpha \geq \alpha_c$  there exists a linear network code such that any data collector can reconstruct. Moreover, randomized network coding at the storage nodes will suffice with high probability.

*Proof:* (sketch) We will show that if  $\alpha < \alpha_c$  the minimum cut from some  $k$  subset of storage nodes to the source  $S$  will be less than  $\mathcal{M}$  and therefore reconstruction will be impossible. In addition when  $\alpha \geq \alpha_c$  the minimum cut will be greater or equal to  $\mathcal{M}$ . Then by Proposition 1 a linear network code exists so that all data collectors can recover. Further randomized network coding will work with probability that can be driven arbitrarily high by increasing the field size.

Therefore it suffices to find the minimum  $\alpha_c$  such that any  $k$  subset of storage nodes has a minimum cut from the source equal to  $\mathcal{M}$ . We proceed via induction on  $n$ , the number of storage nodes. We refer to any subgraph of  $\mathcal{G}$  with  $k$  inputs and  $j \geq k$  outputs as a *box*; a box is called *good* if every  $k$  out of the  $j$  outputs can support an end-to-end flow of  $\mathcal{M}$ . The base case of the induction is trivial if we assume that there are  $k$  storage nodes initially.

For the inductive step, assume we have a good box denoted  $B_{j-1}$  and a newcomer  $X_i$  connects to any  $k$  outputs of  $B_{j-1}$  with edges that have capacity  $\alpha \frac{\mathcal{M}}{k}$ . One needs to show that the new graph with the outputs of  $B_{j-1}$  plus the output of the storage node  $X_i$  will be a good box  $B_j$ . Let  $N(X_i)$  denote the storage nodes where  $X_i$  is connected to. Consider a data collector that connects to  $y_1$  nodes in  $N(X_i)^c$  and  $y_2$  nodes in  $N(X_i)$ , and also to the newcomer (all data collectors that do not connect to the newcomer receive enough flow by the induction hypothesis). We therefore have  $y_1 + y_2 = k - 1$  and also the minimum cut for this data collector is

$$y_1 \alpha \mathcal{M} + y_2 \alpha \mathcal{M} + (k - y_2) \frac{\alpha \mathcal{M}}{k}. \quad (2)$$

To ensure recovery this has to work for every data collector, i.e.

$$y_1 \alpha \mathcal{M} + y_2 \alpha \mathcal{M} + (k - y_2) \frac{\alpha \mathcal{M}}{k} \geq \mathcal{M}, \quad (3)$$

$$\forall y_1, y_2, y_1 + y_2 = k - 1. \quad (4)$$

It is easy to see that  $y_1 = 0$  is the worst case, and from there one obtains that

$$\alpha \geq \frac{1}{k(1 - \frac{1}{k} + \frac{1}{k^2})} =: \alpha_c \quad (5)$$

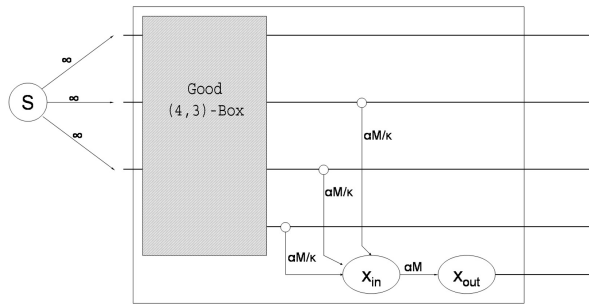
is necessary and sufficient for reconstruction. ■

## V. EVALUATION

In this section, we compare Regenerating Codes with other redundancy management schemes in the context of distributed storage systems. We follow the evaluation methodology of [18], which consists of a simple analytical model whose parameters are obtained from traces of node availability measured in several real distributed systems.

### A. Metrics

Three metrics of importance are *reliability*, *bandwidth*, and *disk usage*. Since bandwidth is generally considered a much more constrained resource than disk space in wide-area environments, we omit an explicit comparison of disk space used by the redundancy management schemes. However,



**Fig. 3.** Illustration of the inductive step. The internal box is good and we want to show that the external box is also good if the newcomer downloads  $1/k\alpha\mathcal{M}$  from the existing nodes the big box is also good.

disk usage would be proportional to bandwidth for all the schemes we evaluate, with the exception of OMMDS, which is the only scheme in which a newcomer stores less on disk than it downloads. We measure reliability in terms of *file availability*, that is, the fraction of time that a file can be reconstructed from the data stored on nodes that are currently available. Another important notion of reliability that we do not evaluate here is *durability*, which measures permanent data loss rate.

### B. Model

We use a model which is intended to capture the average-case bandwidth used to maintain a file in the system, and the resulting average availability of the file. With minor exceptions, this model and the subsequent estimation of its parameters are equivalent to that of [18]. Although this evaluation methodology is a significant simplification of real storage systems, it allows us to compare directly with the conclusions of [18] as well as to calculate precise values for rare events.

The model has two key parameters,  $f$  and  $a$ . First, we assume that in expectation a fraction  $f$  of the nodes storing file data fail per unit time, causing data transfers to repair the lost redundancy. Second, we assume that at any given time while a node is storing data, the node is available with some probability  $a$ . Moreover, the model assumes that the event that a node is available is independent of the availability of all other nodes.

Under these assumptions, we can compute the expected availability and maintenance bandwidth of various redundancy schemes to maintain a file of  $\mathcal{M}$  bytes. We make use of the fact that for all schemes except OMMDS (even Hybrid [18]), the amount of bandwidth used is equal to the amount of redundancy that had to be replaced, which is in expectation  $f$  times the amount of storage used.

**Replication:** If we store  $\mathcal{R}$  replicas of the file, then we store a total of  $\mathcal{R} \cdot \mathcal{M}$  bytes, and in expectation we must replace  $f \cdot \mathcal{R} \cdot \mathcal{M}$  bytes per unit time. The file is unavailable if no replica is available, which happens with probability  $(1 - a)^{\mathcal{R}}$ .

**Ideal Erasure Codes:** For comparison, we show the bandwidth and availability of a hypothetical  $(n, k)$  erasure code strategy which can “magically” create a new packet while transferring just  $\mathcal{M}/k$  bytes (i.e., the size of the packet). Setting  $n = k \cdot \mathcal{R}$ , this strategy sends  $f \cdot \mathcal{R} \cdot \mathcal{M}$  bytes per

Trace	Length (days)	Start date	Mean # nodes up	$f$	$a$
PlanetLab	527	Jan. 2004	303	0.017	0.97
Microsoft	35	Jul. 1999	41970	0.038	0.91
Skype	25	Sep. 2005	710	0.12	0.65
Gnutella	2.5	May 2001	1846	0.30	0.38

TABLE I: The availability traces used in this paper.

unit time and has unavailability probability  $U_{\text{ideal}}(n, k) := \sum_{i=0}^{k-1} \binom{n}{i} a^i (1-a)^{n-i}$ .

**Hybrid:** If we store one full replica plus an  $(n, k)$  erasure code where  $n = k \cdot (\mathcal{R} - 1)$ , then we again store  $\mathcal{R} \cdot \mathcal{M}$  bytes in total, so we transfer  $f \cdot \mathcal{R} \cdot \mathcal{M}$  bytes per unit time in expectation. The file is unavailable if the replica is unavailable and fewer than  $k$  erasure-coded packets are available, which happens with probability  $(1-a) \cdot U_{\text{ideal}}(n, k)$ .

**OMMDS Codes:** A  $(k, n)$  OMMDS Code with redundancy  $\mathcal{R} = n/k$  stores  $\mathcal{R}\mathcal{M}$  bytes in total, so  $f \cdot \mathcal{R} \cdot \mathcal{M}$  bytes must be replaced per unit time. But replacing a fragment requires transferring over the network  $\beta_{\text{OMMDS}} = (n-1)/(n-k)$  times the size of the fragment (see Section III-B), even in the most favorable case when newcomers connect to  $n-1$  nodes to construct a new fragment. This results in  $f \cdot \mathcal{R} \cdot \mathcal{M} \cdot \beta_{\text{OMMDS}}$  bytes sent per unit time, and unavailability  $U_{\text{ideal}}(n, k)$ .

**Regenerating Codes:** A  $(k, n)$  Regenerating Code stores  $\mathcal{M} \cdot n \cdot \beta_{\text{RC}}$  bytes in total (see Section IV). So in expectation  $f \cdot \mathcal{M} \cdot n \cdot \beta_{\text{RC}}$  bytes are transferred per unit time, and the unavailability is again  $U_{\text{ideal}}(n, k)$ .

### C. Estimating $f$ and $a$

In this section we describe how we estimate  $f$ , the fraction of nodes that permanently fail per unit time, and  $a$ , the mean node availability, based on traces of node availability in several distributed systems.

We use four traces of node availability with widely varying characteristics, summarized in Table I, which used periodic network-level probes to determine host availability. The measurements were in four systems representing distinct environments: PlanetLab [20]; a stable, managed network research testbed; desktop PCs at Microsoft Corporation [4]; superpeers in the Skype P2P VoIP network [12], which may approximate the behavior of a set of well-provisioned endhosts, since superpeers are likely selected in part based on available bandwidth [12]; and ordinary peers in the Gnutella filesharing network [19].

It is of key importance for the storage system to distinguish between *transient* failures, in which a node temporarily departs but returns later with its data intact; and *permanent* failures, in which data is lost. Only the latter requires bandwidth-intensive replacement of lost redundancy. Most systems use a *timeout* heuristic: when a node has not responded to network-level probes after some period of time  $t$ , it is considered to have failed permanently. To approximate storage systems' behavior, we use the same heuristic. Node availability  $a$  is then calculated as the mean (over time) fraction of nodes which were available among those which were not considered permanently failed at that time.

The resulting values of  $f$  and  $a$  appear in Table I, where we have fixed the timeout  $t$  at 1 day. Longer timeouts reduce

overall bandwidth costs [18], [5], but begin to impact durability [5] and are more likely to produce artificial effects in the short (2.5-day) Gnutella trace.

## VI. QUANTITATIVE RESULTS AND CONCLUSIONS

Figure 4 shows the tradeoff between mean unavailability and mean maintenance bandwidth in each of the strategies of Section V-B using the values of  $f$  and  $a$  from Section V-C, for  $k = 7$  and  $k = 14$ . Due to space limitations, we omit results for the Microsoft PCs and Skype traces, which lie between those for PlanetLab and Gnutella. Points in the tradeoff space are produced by varying the redundancy factor  $\mathcal{R}$ .

In all cases, OMMDS obtains worse points in the tradeoff space than Hybrid, though it is not much worse for large  $\mathcal{R}$  as shown in the Gnutella results.

Our main conclusion is that our proposed network coding scheme obtains substantial benefits over previous techniques, especially in relatively stable environments. For example, in the PlanetLab trace with  $k = 7$ , RC has about 25% lower bandwidth for the same availability, or more than 3 orders of magnitude lower unavailability with the same bandwidth. The difference is even greater for  $k = 14$ .

RC's reduction in bandwidth compared with Hybrid diminishes as the environment becomes less stable; in the most extreme case of the Gnutella trace, RC can actually be very slightly *worse*. The reason can be seen by comparing the two schemes with Ideal Erasure Codes. For fixed  $k$  and  $n$ , both RC and Hybrid have roughly the same availability (Hybrid is slightly better due to the extra replica). However, in terms of bandwidth as we scale  $n$ , RC has a small *constant factor* overhead compared with Ideal Erasure codes, while Hybrid has a rather large but only *additive* overhead due to the single extra replica. For large enough  $n$ , such as is necessary in Gnutella, the additive overhead wins out. But such scenarios are unlikely to be practical in any case due to the high bandwidth required of all schemes.

However, RC still has some drawbacks. First, constructing a new packet, or reconstructing the entire file, requires communication with  $k$  nodes rather than one (in Hybrid, the node holding the single replica). This adds overhead that could be significant for sufficiently small files or sufficiently large  $k$ . Perhaps more importantly, there is a slight increase in total data transferred to *read* the file, roughly 14% for  $k = 7$  but diminishing to 7.1% for  $k = 14$  and 3.1% for  $k = 32$ . Thus, if the frequency that a file is read is sufficiently high and  $k$  is sufficiently small, this inefficiency could overwhelm the reduction in maintenance bandwidth.

If the target application is archival storage or backup, files are likely to be large and infrequently read. We believe this is one case in which RC is likely to be a significant win over both Hybrid and replication.

Our results suggest that network coding can provide a significant reduction in maintenance bandwidth and also simplify system architecture since only one type of redundancy needs to be maintained. This addresses the two principal disadvantages of using erasure coding discussed in [18], and therefore we believe that regenerating codes are a promising design choice for certain peer-to-peer storage systems.

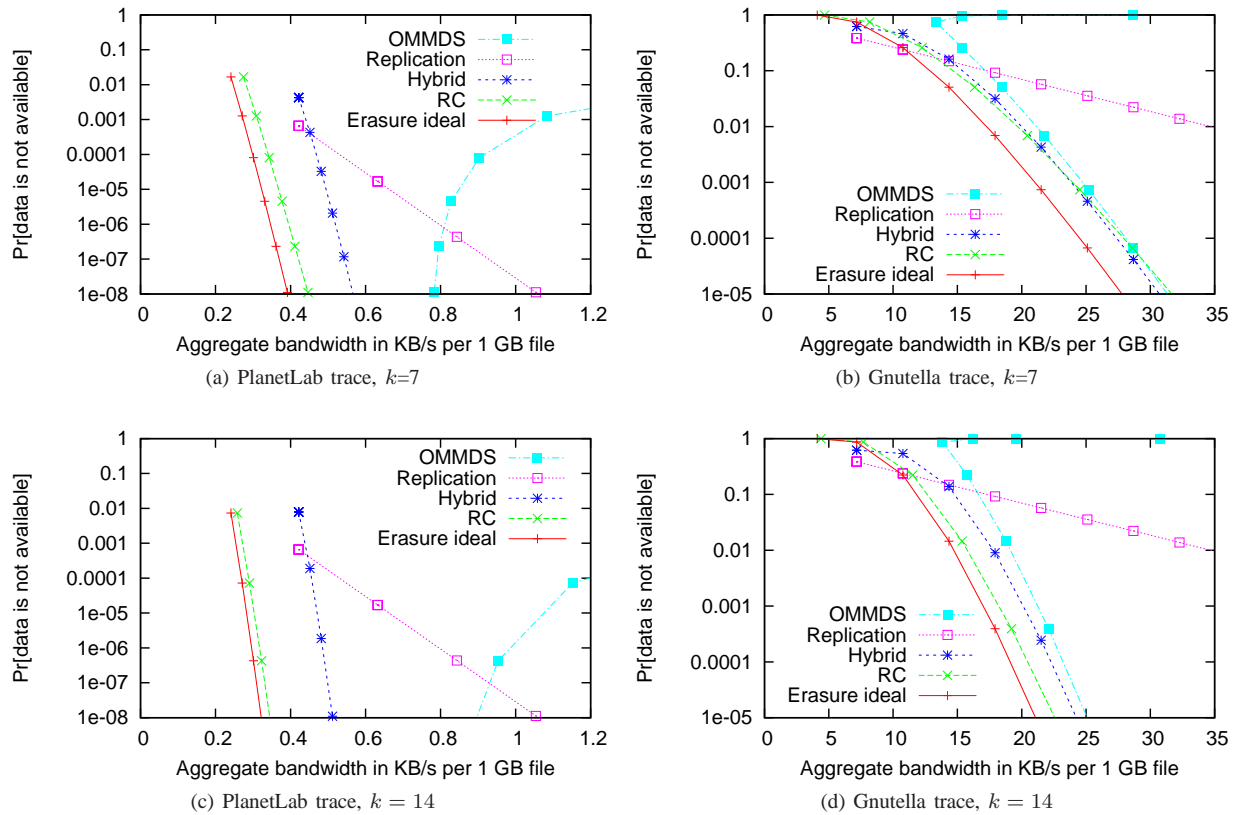


Fig. 4: Availability-bandwidth tradeoff produced by varying redundancy  $\mathcal{R}$ , with parameters derived from the traces.

#### REFERENCES

- [1] S. Acedanski, S. Deb, M. Médard, and R. Koetter. How good is random linear coding based distributed networked storage. In *NetCod*, 2005.
- [2] R. Ahlswede, N. Cai, S.-Y. R. Li, and R. W. Yeung. Network information flow. *IEEE Trans. Info. Theory*, 46(4):1204–1216, July 2000.
- [3] Ranjita Bhagwan, Kiran Tati, Yu-Chung Cheng, Stefan Savage, and Geoffrey M. Voelker. Total recall: System support for automated availability management. In *NSDI*, 2004.
- [4] William J. Bolosky, John R. Douceur, David Ely, and Marvin Theimer. Feasibility of a serverless distributed file system deployed on an existing set of desktop PCs. In *Proc. SIGMETRICS*, 2000.
- [5] Byung-Gon Chun, Frank Dabek, Andreas Haeberlen, Emil Sit, Hakim Weatherspoon, M. Frans Kaashoek, John Kubiatowicz, and Robert Morris. Efficient replica maintenance for distributed storage systems. In *NSDI*, 2006.
- [6] F. Dabek, J. Li, E. Sit, J. Robertson, M. Kaashoek, and R. Morris. Designing a dht for low latency and high throughput, 2004.
- [7] A. G. Dimakis, V. Prabhakaran, and K. Ramchandran. Ubiquitous Access to Distributed Data in Large-Scale Sensor Networks through Decentralized Erasure Codes. In *Proc. IEEE/ACM Int. Symposium on Information Processing in Sensor Networks (IPSN)*, April 2005.
- [8] A. G. Dimakis, V. Prabhakaran, and K. Ramchandran. Decentralized erasure codes for distributed networked storage. In *IEEE Trans on Information Theory*, June 2006.
- [9] A.G. Dimakis, P.B. Godfrey, M.J. Wainwright, and K. Ramchandran. Network coding for distributed storage systems. In *Proceedings of IEEE INFOCOM (to appear)*, 2007.
- [10] C. Fragouli, J.Y. Le Boudec, and J. Widmer. On the benefits of network coding for wireless applications. *NetCod*, 2006.
- [11] C. Gkantsidis and P. Rodriguez. Network coding for large scale content distribution. *Proceedings of IEEE Infocom*, 2005.
- [12] Saikat Guha, Neil Daswani, and Ravi Jain. An experimental study of the Skype peer-to-peer VoIP system. In *IPTPS*, 2006.
- [13] T. Ho, M. Médard, R. Koetter, D. Karger, M. Effros, J. Shi, and B. Leong. A random linear network coding approach to multicast. *Submitted for publication, IEEE Trans. Info. Theory*, 2006.
- [14] A. Kamra, J. Feldman, V. Misra, and D. Rubenstein. Growth codes: Maximizing sensor network data persistence. *ACM SIGCOMM*, 2006.
- [15] S.-Y. R. Li, R. W. Yeung, and N. Cai. Linear network coding. *IEEE Trans. on Information Theory*, 49:371–381, February 2003.
- [16] S. Rhea, P. Eaton, D. Geels, H. Weatherspoon, B. Zhao, and J. Kubiatowicz. Pond: the OceanStore prototype. In *Proc. USENIX File and Storage Technologies (FAST)*, 2003.
- [17] S. Rhea, C. Wells, P. Eaton, D. Geels, B. Zhao, H. Weatherspoon, and J. Kubiatowicz. Maintenance-free global data storage. *IEEE Internet Computing*, pages 40–49, September 2001.
- [18] R. Rodrigues and B. Liskov. High availability in DHTs: Erasure coding vs. replication. In *Proc. IPTPS*, 2005.
- [19] Stefan Saroiu, P. Krishna Gummadi, and Steven D. Gribble. A Measurement Study of Peer-to-Peer File Sharing Systems. In *Proc. MMCN*, San Jose, CA, USA, January 2002.
- [20] Jeremy Stribling. Planetlab all pairs ping. <http://infospect.planetlab.org/pings>.
- [21] D. Wang, Q. Zhang, and J. Liu. Partial network coding: Theory and application for continuous sensor data collection. *Fourteenth IEEE International Workshop on Quality of Service (IWQoS)*, 2006.
- [22] Hakim Weatherspoon and John D. Kubiatowicz. Erasure coding vs. replication: a quantitative comparison. In *Proc. IPTPS*, 2002.
- [23] C. Wu and B. Li. Echelon: Peer-to-peer network diagnosis with network coding. *Fourteenth IEEE International Workshop on Quality of Service (IWQoS)*, 2006.
- [24] X. Zhang, G. Neglia, J. Kurose, and D. Towsley. On the benefits of random linear coding for unicast applications in disruption tolerant networks. *Second Workshop on Network Coding, Theory, and Applications (NETCOD)*, 2006.